



UANL

UNIVERSIDAD AUTÓNOMA DE NUEVO LEÓN

FCFM



FACULTAD DE CIENCIAS FÍSICO MATEMÁTICAS

Universidad Autónoma de Nuevo León

Facultad de Ciencias Físico Matemáticas

Producto Integrador de Aprendizaje

Aprendizaje automático

Nombre: Sergio Andrés Elizondo Rodríguez

Grupo: 001

Maestro: José Anastacio Hernández Saldaña

Viernes 26 de Julio del 2024

Reporte de resultados

Objetivo: Encontrar el mejor modelo de clasificación utilizando validación cruzada y el criterio de área bajo la curva característica operativa del receptor (*ROC AUC*), logrando un valor mayor a 0.75 en el conjunto de validación y prueba.

Datos utilizados:

- **Conjunto de entrenamiento:** 9239 observaciones, 10 columnas.
- **Conjunto de prueba:** 2309 observaciones, 9 columnas.

Procesamiento previo de los datos

- Se convirtió la variable objetivo (*engagement*) a valores numéricos.
- Se separaron las características (*X*) y las etiquetas (*y*) en el conjunto de entrenamiento.
- Se normalizaron las características utilizando un escalador (*standard scaler*).
- Se dividió el conjunto etiquetado (*entrenamiento*) en datos de entrenamiento y de validación (80% y 20%, respectivamente).

Selección de modelo

Se evaluaron los siguientes modelos:

- Regresión logística
- *K* vecinos más cercanos (*KNN*)
- Clasificador de margen máximo (*support vector classifier*)
- Árbol de decisión
- Bosque aleatorio (*random forest*)
- Potencialización de gradiente (*gradient boosting*)
- *XGBoost*
- *LightGBM*

Para la evaluación, se definieron espacios de búsqueda de hiperparámetros específicos para cada modelo, optimizando su rendimiento mediante validación cruzada de 5 pliegues en una búsqueda en cuadrícula (*grid search cross-validation*) con el criterio *ROC AUC*.

Resultados

A continuación se presentan los resultados de la métrica *ROC AUC* obtenidos en el conjunto de validación para los mejores hiperparámetros de cada modelo:

| Modelo | Mejores hiperparámetros | Área bajo la curva operativa del receptor (<i>ROC AUC</i>) |
|--|--|--|
| Árbol de decisión | <ul style="list-style-type: none"> • Ajuste de pesos: Ninguno • Criterio de división: Entropía • Profundidad máxima: Ilimitada • Máximo de características por división: 80% • Máximo de nodos hoja: 40 • Umbral para reducción de impureza: 0 • Mínimo de muestras por hoja: 5 • Mínimo de muestras por división: 5 • Fracción mínima del peso por hoja: 0.01 • Estrategia de división: Elegir la mejor | 85.93% |
| <i>K</i> vecinos más cercanos (<i>KNN</i>) | <ul style="list-style-type: none"> • Vecinos a considerar: 55 • Potencia de la distancia de Minkowski: 2 (distancia euclidiana) • Ponderación: Por distancia | 87.68% |
| Regresión logística | <ul style="list-style-type: none"> • Parámetro de regularización inversa: 0.02 • Ajuste de pesos: Balanceado • Tipo de regularización: L2 (<i>ridge</i>) • Algoritmo de optimización: <i>liblinear</i> | 87.84% |
| Clasificador de margen máximo (<i>support vector classifier</i>) | <ul style="list-style-type: none"> • Parámetro de regularización: 3 • Núcleo (<i>kernel</i>): Polinómico • Constante del núcleo: 19 • Grado del polinomio: 3 • Coefficiente γ: 0.04 | 89.07% |
| Bosque aleatorio (<i>random forest</i>) | <ul style="list-style-type: none"> • Ajuste de pesos: Ninguno • Criterio de división: Entropía • Profundidad máxima: 10 • Máximo de características por división: 75% • Mínimo de muestras por hoja: 1 • Mínimo de muestras por división: 2 • Número de árboles: 200 | 89.25% |
| <i>XGBoost</i> | <ul style="list-style-type: none"> • Amplificador (<i>booster</i>): Regresión múltiple aditiva (<i>DART</i>) • Reducción mínima (γ): 0.1 • Tasa de aprendizaje: 0.05 • Profundidad máxima: 5 • Peso mínimo por hoja: 2 • Restricciones monotónicas: (1, 0) • Número de árboles: 97 | 90.31% |

| | | |
|--|--|--------|
| | <ul style="list-style-type: none"> • Función de pérdida: Logística binaria • Término de regularización L1 (LASSO): 1 • Término de regularización L2 (ridge): 0 • Peso de las clases positivas: 2 • Fracción de muestras para entrenamiento: 90% | |
| <i>LightGBM</i> | <ul style="list-style-type: none"> • Tasa de aprendizaje: 0.01 • Profundidad máxima: Ilimitada • Mínimo de muestras por hoja: 30 • Número de árboles: 350 • Máximo de hojas por árbol: 40 | 90.31% |
| Potencialización de gradiente (<i>gradient boosting</i>) | <ul style="list-style-type: none"> • Tasa de aprendizaje: 0.1 • Función de pérdida: Exponencial • Profundidad máxima: 5 • Mínimo de muestras por hoja: 1 • Mínimo de muestras por división: 2 • Fracción mínima del peso por hoja: 0.03 • Número de árboles: 100 | 90.45% |

El modelo de **potencialización de gradiente (*gradient boosting*)** fue el que obtuvo el mejor rendimiento ($ROC\ AUC = 90.45\%$).

Evaluación

Para la evaluación final, se utilizó el mejor modelo encontrado (*gradient boosting*) para predecir las probabilidades en el conjunto de prueba. Las predicciones se guardaron en el archivo “resultados.csv”, asignando la probabilidad de participación (*engagement*) predicha para cada identificador (*id*).

Conclusión

El modelo de potencialización de gradiente (*gradient boosting*) demostró ser el más efectivo en la evaluación de área bajo la curva operativa del receptor ($ROC\ AUC$), con una métrica de 0.9045, superando el umbral de 0.75 en el conjunto de validación, lo cual confirma que el modelo seleccionado es adecuado para la tarea de clasificación planteada. Este rendimiento superior puede atribuirse a varios factores, como la optimización de una función de pérdida robusta frente a valores atípicos y la consideración de un gran número de árboles de decisión, lo cual permite que los errores de un árbol se minimicen por los demás árboles, permitiendo a su vez capturar relaciones no lineales entre los datos.

Además el ajuste fino de hiperparámetros permitió mejorar la capacidad del modelo de capturar patrones específicos de nuestro conjunto de datos sin caer en el sobreajuste (gracias a la validación cruzada). Es necesario contar con las etiquetas en el conjunto de prueba para evaluar si efectivamente el buen comportamiento del modelo se sigue replicando con los mismos.